



Computer Science and Artificial Intelligence Laboratory

Technical Report

MIT-CSAIL-TR-2007-051

November 1, 2007

Towards Feature Selection In Actor-Critic Algorithms

Khashayar Rohanimanesh, Nicholas Roy, and Russ Tedrake

Report Documentation Page			Form Approved OMB No. 0704-0188		
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE 01 NOV 2007		2. REPORT TYPE		3. DATES COVERED 00-00-2007 to 00-00-2007	
4. TITLE AND SUBTITLE Towards Feature Selection In Actor-Critic Algorithms			5a. CONTRACT NUMBER		
			5b. GRANT NUMBER		
			5c. PROGRAM ELEMENT NUMBER		
6. AUTHOR(S)			5d. PROJECT NUMBER		
			5e. TASK NUMBER		
			5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Massachusetts Institute of Technology, Computer Science and Artificial Intelligence Laboratory ,32 Vassar Street,Cambridge,MA,02139			8. PERFORMING ORGANIZATION REPORT NUMBER		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)			10. SPONSOR/MONITOR'S ACRONYM(S)		
			11. SPONSOR/MONITOR'S REPORT NUMBER(S)		
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT Choosing features for the critic in actor-critic algorithms with function approximation is known to be a challenge. Too few critic features can lead to degeneracy of the actor gradient, and too many features may lead to slower convergence of the learner. In this paper, we show that a well-studied class of actor policies satisfy the known requirements for convergence when the actor features are selected carefully. We demonstrate that two popular representations for value methods - the barycentric interpolators and the graph Laplacian proto-value functions - can be used to represent the actor in order to satisfy these conditions. A consequence of this work is a generalization of the proto-value function methods to the continuous action actor-critic domain. Finally, we analyze the performance of this approach using a simulation of a torque-limited inverted pendulum.					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT Same as Report (SAR)	18. NUMBER OF PAGES 11	19a. NAME OF RESPONSIBLE PERSON
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified			

Towards Feature Selection In Actor-Critic Algorithms

Khashayar Rohanimanesh, Nicholas Roy, Russ Tedrake
{khash,nickroy,russt}@mit.edu
Massachusetts Institute of Technology

November 1, 2007

Abstract

Choosing features for the critic in actor-critic algorithms with function approximation is known to be a challenge. Too few critic features can lead to degeneracy of the actor gradient, and too many features may lead to slower convergence of the learner. In this paper, we show that a well-studied class of actor policies satisfy the known requirements for convergence when the actor features are selected carefully. We demonstrate that two popular representations for value methods - the barycentric interpolators and the graph Laplacian proto-value functions - can be used to represent the actor in order to satisfy these conditions. A consequence of this work is a generalization of the proto-value function methods to the continuous action actor-critic domain. Finally, we analyze the performance of this approach using a simulation of a torque-limited inverted pendulum.

1 Introduction

Actor-Critic (AC) algorithms, initially proposed by Barto et al. (1983), aim at combining the strong elements of the two major classes of reinforcement learning algorithms – namely the value-based methods and the policy search methods. As in value-based methods, the critic component maintains a value function, and as in policy search methods, the actor component maintains a separate parameterized stochastic policy from which the actions are drawn. This combination may offer the convergence guarantees which are characteristic of the policy gradient algorithms as well as an improved convergence rate because the critic can be used to reduce the variance of the policy update (Konda and Tsitsiklis, 2003).

Recent AC algorithms use a function approximation architecture to maintain both the actor policy and the critic (state-action) value function, relying on Temporal Difference (TD) learning methods to update the critic parameters. Konda and Tsitsiklis (2000) and Sutton et al. (2000) showed that in order to compute the gradient of the performance function (typically using the average cost criterion) with respect to the parameters of a stochastic policy $\mu_\theta(\mathbf{x}, \mathbf{u})$ it suffices to compute the projection of the state-action value function onto a sub-space Ψ spanned by the vectors $\psi_\theta^i(x, u) = \frac{\partial}{\partial \theta_i} \log \mu_\theta(\mathbf{x}, \mathbf{u})$. Konda and Tsitsiklis (2003) also noted that for certain values of the policy parameters θ , it is possible that the vectors ψ_θ^i are either close to zero, or almost linearly dependent. In these situations the projection onto Ψ becomes ill-conditioned, providing no useful gradient information, and the algorithm can become unstable. As a remedy for this problem the authors suggested the use of a richer, higher dimensional set of critic features which contain the space Ψ as a proper subset.

In this paper, we attempt to design features which span Ψ and preserve linear independence without increasing the dimensionality of the critic. In particular, we investigate stochastic actor policies represented by a family of Gaussian distributions where the mean of the distribution is linearly parameterized using a set of a fixed basis functions. For this parameterization, we show that if the basis functions in the actor are selected to be linearly independent, then the minimal set of critic features which naturally satisfy the containment condition also form a linearly independent basis set. Additionally, if the actor basis set is linearly independent of the function $\mathbf{1}$, then the critic features satisfy a weak version of the non-zero projection condition specified by Konda and Tsitsiklis (2003). This suggests that feature sets which have been proposed for representing value functions, such as the *proto-value-functions* (Mahadevan and Maggioni, 2006), may also have promise as features for actor-critic algorithms. This extends the proto-value function approach, which traditionally works by discretizing the action space, to a continuous action actor-critic domain.

The rest of this paper is organized as follows: in Section 2 we provide a brief review of the AC algorithms with function approximation. In Section 3 we present our main theoretical results by investigating a family of parameterized Gaussian policies. In Section 4 we consider candidate features which satisfy these results. In Section 5 we describe an empirical evaluation of our approach in a simulated control domain. Finally, we discuss some implications and future directions in Section 6 and conclude in Section 7.

2 Preliminaries

In this section we present a brief overview of the AC algorithms with function approximation adapted from Konda and Tsitsiklis (2003). Assume that the problem is modeled as a Markov decision process $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{P}, \mathcal{C} \rangle$, where \mathcal{X} is the state space, \mathcal{U} is the action space, $\mathcal{P}(x'|x, u)$ is the transition probability function, $\mathcal{C} : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ is the one step cost function, and μ_θ is a stochastic policy parameterized by $\theta \in \mathbb{R}^n$ where $\mu_\theta(u|x)$ gives the probability of selecting an action u in state x , parameterized by the vector $\theta \in \mathbb{R}^n$. We also assume that for every $\theta \in \mathbb{R}^n$, the Markov chains $\{X_k\}$ and $\{X_k, U_k\}$ are irreducible and aperiodic, with stationary probabilities $\pi_\theta(x)$ and $\eta_\theta(x, u) = \pi_\theta(x)\mu_\theta(u|x)$ respectively.

The average cost function $\bar{\alpha}(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}$ can be defined:

$$\bar{\alpha}(\theta) = \sum_{x \in \mathcal{X}, u \in \mathcal{U}} c(x, u) \eta_\theta(x, u)$$

For each $\theta \in \mathbb{R}^n$, let $\mathcal{V}_\theta : \mathcal{X} \rightarrow \mathbb{R}$, and $\mathcal{Q}_\theta : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ be the differential state, and the differential state-action cost functions that are solution to the corresponding Poisson equations in a standard average cost setting. Then, following the results of Marbach and Tsitsiklis (1998), the gradient of the average cost function can be expressed as:

$$\nabla_\theta \bar{\alpha}(\theta) = \sum_{x, u} \eta_\theta(x, u) \mathcal{Q}_\theta(x, u) \psi_\theta(x, u) \quad (1)$$

where:

$$\psi_\theta(x, u) = \nabla_\theta \ln \mu_\theta(u|x) \quad (2)$$

The i th component of $\psi_\theta, \psi_\theta^i(x, u)$ is the one-step *eligibility* of parameter i in state-action pair x and u :

$$\psi_\theta^i(x, u) = \frac{\partial}{\partial \theta_i} \ln \mu_\theta(u|x). \quad (3)$$

Following Konda and Tsitsiklis (2003), we will assume that \mathcal{X} and \mathcal{U} are discrete, countably infinite sets. We will therefore refer to ψ_θ^i as the *actor eligibility vector*, a vector in $\mathbb{R}^{|\mathcal{X}||\mathcal{U}|}$. For any $\theta \in \mathbb{R}^n$, the inner product $\langle \cdot, \cdot \rangle_\theta$ of two real-valued functions $\mathcal{Q}_1, \mathcal{Q}_2$ on $\mathcal{X} \times \mathcal{U}$, also viewed as vectors in $\mathbb{R}^{|\mathcal{X}||\mathcal{U}|}$, can be defined by:

$$\langle \mathcal{Q}_1, \mathcal{Q}_2 \rangle_\theta = \sum_{x, u} \eta_\theta(x, u) \mathcal{Q}_1(x, u) \mathcal{Q}_2(x, u)$$

and let $\| \cdot \|_\theta$ denote the norm induced by this inner product on $\mathbb{R}^{|\mathcal{X}||\mathcal{U}|}$. Now, we can rewrite Equation 1 as:

$$\frac{\partial}{\partial \theta_i} \bar{\alpha}(\theta) = \langle \mathcal{Q}_\theta, \psi_\theta^i \rangle_\theta, \quad i = 1, \dots, n.$$

For each $\theta \in \mathbb{R}^n$, let Ψ_θ denote the span of the vectors $\{\psi_\theta^i; 1 \leq i \leq n\}$ in $\mathbb{R}^{|\mathcal{X}||\mathcal{U}|}$. An important observation is that although the gradient of $\bar{\alpha}$ depends on the function \mathcal{Q}_θ , which is a vector in a possibly very high-dimensional space $\mathbb{R}^{|\mathcal{X}||\mathcal{U}|}$, the dependence is only through its inner products with vectors in Ψ_θ . Thus, instead of “learning” the function \mathcal{Q}_θ , it suffices to learn its projection on the low-dimensional sub-space Ψ_θ .

Konda and Tsitsiklis (2003) consider actor-critic algorithms where the critic is a TD algorithm with a linearly parameterized approximation architecture for the Q -value function that admits the linear-additive form:

$$Q_\theta^r(x, u) = \sum_{j=1}^m r^j \phi_\theta^j(x, u) \quad (4)$$

where $r = (r^1, \dots, r^m) \in \mathbb{R}^m$ is the parameter vector of the critic. The critic features $\phi_\theta^j, j = 1, \dots, m$ depend on the actor parameter vector and are chosen so that the following assumptions are satisfied: (1) For every $(x, u) \in \mathcal{X} \times \mathcal{U}$, the map $\theta \rightarrow \phi_\theta(x, u)$ is bounded and differentiable; (2) The span of the vectors $\phi_\theta^j (j = 1, \dots, m)$ in $\mathbb{R}^{|\mathcal{X}||\mathcal{U}|}$ denoted by Φ_θ , contains Ψ_θ .

As noted by Konda and Tsitsiklis (2003), one trivial choice for satisfying the second condition would be to set $\Psi = \Phi$, or in other words to set critic features as $\phi_\theta^i = \psi_\theta^i$. However, it is possible that for some values of θ , the features ψ_θ^i are either close to zero, or almost linearly dependent. In these situations the projection of Q_θ^r onto Ψ becomes ill-conditioned, providing no useful gradient information, and therefore the algorithm may become unstable. Konda and Tsitsiklis (2003) suggest some ideas to remedy to this problem. In particular, the troublesome situations are avoided if the following condition is satisfied: (3) There exists $a > 0$, such that for every $r \in \mathbb{R}^m$, and $\theta \in \mathbb{R}^n$:

$$\| r' \hat{\phi}_\theta \|^2 \geq a |r|^2$$

where $\hat{\phi} = \{\hat{\phi}^i\}_{i=1}^m$ are defined as:

$$\hat{\phi}_\theta^i(x, u) = \phi_\theta^i(x, u) - \sum_{\bar{x}, \bar{u}} \eta_\theta(\bar{x}, \bar{u}) \phi_\theta^i(\bar{x}, \bar{u}) \quad (5)$$

This condition can be roughly explained as follows: the new functions $\hat{\phi}_\theta^i$ can be viewed as the original critic features with their expected value (with respect to the distribution $\eta_\theta(x, u)$) removed. In order to ensure that the projection of Q_θ^r onto Ψ contains some gradient information for the actor (and to avoid instability), the set $\hat{\phi}_\theta$ must be uniformly bounded away from zero. Given these conditions, Konda and Tsitsiklis (2003) prove convergence for of the most common form for the actor-critic update (see Konda and Tsitsiklis (2003, p.1148) for the updates).

Konda and Tsitsiklis (2003) go on to propose adding additional features to the critic as a remedy, but satisfying this condition is still a difficult problem. To the best of our knowledge there is no general systematic approach for choosing a set of critic features that satisfies this third condition. In the next section, we will address this issue for one commonly used policy class.

3 Our Approach

We consider the following popular Gaussian probabilistic policy structure parameterized by θ :

$$\mu_\theta(\mathbf{u}|\mathbf{x}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(\mathbf{u} - \mathbf{m}_\theta(\mathbf{x}))^T \Sigma^{-1} (\mathbf{u} - \mathbf{m}_\theta(\mathbf{x}))\right\} \quad (6)$$

where $\mathbf{u} \in \mathbb{R}^k$ is a multi-dimensional action vector. The vector $\mathbf{m}_\theta(\mathbf{x}) \in \mathbb{R}^k$ is the mean of the distribution that is parameterized by θ :

$$\mathbf{m}_\theta^i(\mathbf{x}) = \sum_{j=1}^n \theta^{ij} \rho^j(\mathbf{x}), \quad i = 1, \dots, k$$

where in this setting $\theta \in \mathbb{R}^{k \times n}$. The functions $\rho^j(\mathbf{x})$ are a set of actor features defined over the states. For simplicity, in this paper we only investigate the case where $\Sigma = \sigma_0^2 \mathbf{I}$. In this case Equation 6 simplifies to:

$$\mu_\theta(\mathbf{u}|\mathbf{x}) = \frac{1}{(2\pi)^{\frac{k}{2}} \sigma_0^k} \exp\left\{-\frac{1}{2\sigma_0^2}(\mathbf{u} - \mathbf{m}_\theta(\mathbf{x}))^T (\mathbf{u} - \mathbf{m}_\theta(\mathbf{x}))\right\}$$

Using Equation 3 we can compute the actor eligibility vectors as follows:

$$\begin{aligned}
\psi_{\theta}^{ij}(\mathbf{x}, \mathbf{u}) &= \frac{\partial}{\partial \theta^{ij}} \ln \mu_{\theta}(\mathbf{u}|\mathbf{x}) \\
&= \frac{\partial}{\partial \theta^{ij}} \left[-\ln((2\pi)^{\frac{k}{2}} \sigma_0^k) - \frac{1}{2\sigma_0^2} (\mathbf{u} - \mathbf{m}_{\theta}(\mathbf{x}))^T (\mathbf{u} - \mathbf{m}_{\theta}(\mathbf{x})) \right] \\
&= \frac{1}{\sigma_0^2} (\mathbf{u} - \mathbf{m}_{\theta}(\mathbf{x}))^T \frac{\partial}{\partial \theta^{ij}} \mathbf{m}_{\theta}(\mathbf{x}) \\
&= \frac{1}{\sigma_0^2} (u_i - \mathbf{m}_{\theta}^i(\mathbf{x})) \rho^j(\mathbf{x}) \\
&= \kappa_{\theta}^i(\mathbf{x}, \mathbf{u}) \rho^j(\mathbf{x})
\end{aligned} \tag{7}$$

where $\kappa_{\theta}^i(\mathbf{x}, \mathbf{u}) = \frac{(u_i - \mathbf{m}_{\theta}^i(\mathbf{x}))}{\sigma_0^2}$. In order to satisfy the condition (2) in the previous section (Φ should properly contain Ψ), we apply the straightforward solution of setting $\phi^{ij} = \psi^{ij}$ for $i = 1, \dots, k$ and $j = 1, \dots, n$. This selection also guarantees that the mapping from θ to ϕ_{θ} is bounded and differentiable, from condition (1). In Proposition 1, we show that for the particular choice of policy structure that we have chosen, if the actor features, $\rho^j(\mathbf{x})$, are linearly independent, then the critic features will also be linearly independent.

Proposition 1: If the functions $\rho = \{\rho^j\}_{j=1}^n$ are linearly independent, then the set of critic feature functions ϕ^{ij} will form a linearly independent set of functions.

Proof: We prove by contradiction that if the above condition holds, then the set of actor eligibility functions ψ^{ij} (and therefore also ϕ^{ij}) are linearly independent. Assume that the functions ψ^{ij} are linearly dependent. Then there exists $\alpha = \{\alpha_{ij} \in \mathbb{R}\}_{i=1, j=1}^{k, n}$ such that

$$\sum_{i=1, j=1}^{k, n} \alpha_{ij} \psi^{ij}(\mathbf{x}, \mathbf{u}) = 0, \quad \forall \mathbf{x} \in \mathcal{X}, \forall \mathbf{u} \in \mathcal{U},$$

and $\|\alpha\|^2 > 0$. Substituting the right hand side of the Equation 7 for $\psi^{ij}(\mathbf{x}, \mathbf{u})$ yields:

$$\sum_{i=1, j=1}^{k, n} \alpha_{ij} \kappa_{\theta}^i(\mathbf{x}, \mathbf{u}) \rho^j(\mathbf{x}) = 0, \quad \forall \mathbf{x} \in \mathcal{X}, \forall \mathbf{u} \in \mathcal{U}.$$

By regrouping terms we obtain:

$$\sum_{j=1}^n \left(\sum_{i=1}^k \alpha_{ij} \kappa_{\theta}^i(\mathbf{x}, \mathbf{u}) \right) \rho^j(\mathbf{x}) = 0, \quad \forall \mathbf{x} \in \mathcal{X}, \forall \mathbf{u} \in \mathcal{U}.$$

Since according to the assumption the functions $\rho = \{\rho^j\}_{j=1}^n$ are linearly independent, then the following condition must hold:

$$\sum_{i=1}^k \alpha_{ij} \kappa_{\theta}^i(\mathbf{x}, \mathbf{u}) = 0, \quad \forall j = 1, \dots, n, \forall \mathbf{x} \in \mathcal{X}, \forall \mathbf{u} \in \mathcal{U} \tag{8}$$

But for every i , there exists an (\mathbf{x}, \mathbf{u}) such that $\kappa_{\theta}^i(\mathbf{x}, \mathbf{u}) \neq 0$:

$$\kappa_{\theta}^i(\mathbf{x}, \mathbf{m}_{\theta}(\mathbf{x}) + \epsilon_i \mathbf{1}) = \frac{\epsilon_i}{\sigma_0^2}, \tag{9}$$

where $\mathbf{1}$ is the $k \times 1$ vector of ones, and $\epsilon_i \neq 0$. Note that the above condition holds for all $\epsilon_i \in \mathbb{R} - \{0\}$. Now, define a $(k \times 1)$ vector \mathbf{h}_l (for $l = 1, \dots, k$) as:

$$\mathbf{h}_l(j) = \begin{cases} \epsilon & \text{if } j \neq l \\ 2\epsilon & \text{if } j = l \end{cases} \tag{10}$$

for some $\epsilon > 0$. Based on Equation 9, if we choose $\mathbf{u} = \mathbf{m}_\theta(\mathbf{x}) + \mathbf{h}_l$ in Equation 8, we obtain:

$$\frac{1}{\sigma_0^2} \mathbf{h}_l^T \bar{\alpha}_{ij} = 0, \quad \forall l = 1, \dots, k$$

where $\bar{\alpha}_{ij} = [\alpha_{1j}, \alpha_{2j}, \dots, \alpha_{kj}]^T$. This gives us the following system of equations (for a fixed value of j):

$$\mathbf{A} \bar{\alpha}_{ij} = \mathbf{0}, \quad i = 1, \dots, k \quad (11)$$

where $\mathbf{A}_{k \times k} = [\mathbf{h}_1 \mathbf{h}_2 \dots \mathbf{h}_k]^T$. Note that for the particular choice of the vectors \mathbf{h}_l (Equation 10), the matrix \mathbf{A} has a full-rank (since the vectors \mathbf{h}_l are linearly independent), and thus the only solution to the Equation 11 is $\bar{\alpha}_{ij} = \mathbf{0}$. This means that $\alpha_{ij} = 0$ (for all i, j), and thus $\|\alpha\|^2 = 0$. By contradiction, ψ^{ij} (and therefore ϕ^{ij}) must be linearly independent.

Proposition 1 provides a mechanism for ensuring that the θ -dependent critic features remain linearly independent for all θ 's, thereby avoiding a major source of potential instabilities in the AC algorithm. However, to meet the strict conditions from Konda and Tsitsiklis (2003), we should also demonstrate that the critic features are uniformly bounded away from zero. Proposition 2 allows us to demonstrate that a set of actor features that is also linearly independent with the function $\underline{1}$ satisfies the weak form of condition (3).

Proposition 2: If the functions $\underline{1}$ and $\rho = \{\rho^j\}_{j=1}^n$, $j = 1, \dots, n$ are linearly independent, then the set of critic feature functions ϕ^{ij} and the function $\underline{1}$, will also form a linearly independent set of functions.

Proof (sketch): We follow the proof of the proposition 1. Assume that the functions ψ^{ij} are linearly dependent. Then there exists $\alpha = \{\alpha_{ij} \in \mathbb{R}\}_{i=1, j=1}^{k, n} \cup \{\alpha_1 \in \mathbb{R}\}$ such that:

$$\sum_{i=1, j=1}^{k, n} \alpha_{ij} \psi^{ij}(\mathbf{x}, \mathbf{u}) + \alpha_1 \underline{1} = 0, \quad \forall \mathbf{x} \in \mathcal{X}, \forall \mathbf{u} \in \mathcal{U},$$

and $\|\alpha\|^2 > 0$. Following the same steps as in proof of proposition 1, we obtain:

$$\sum_{j=1}^n \left(\sum_{i=1}^k \alpha_{ij} \kappa_\theta^i(\mathbf{x}, \mathbf{u}) \right) \rho^j(\mathbf{x}) + \alpha_1 \underline{1} = 0, \quad \forall \mathbf{x} \in \mathcal{X}, \forall \mathbf{u} \in \mathcal{U}. \quad (12)$$

Since according to the assumption the functions $\rho = \{\rho^j\}_{j=1}^n$ and $\underline{1}$ are linearly independent, then $\alpha_1 = 0$. Following the rest of the steps in proof of the proposition 1, it can be also established that $\alpha_{ij} = 0$ (for all i, j). This completes the proof.

Konda and Tsitsiklis (2003) prove that if the functions $\underline{1}$ and the critic features ϕ_θ^i are linearly independent for each θ , then there exists a positive function $a(\theta)$ such that:

$$\|r' \hat{\phi}_\theta\|_\theta^2 \geq a(\theta) |r|^2 \quad (13)$$

(refer to Section 2, Equation 5 for the definition of $\hat{\phi}_\theta$). This is the weak form of the non-zero projection property.

Finally, it should be noted that it is also possible to tune the standard-deviation of the policy distribution, σ_0 , as a function of state using additional policy parameters, \mathbf{w} . If we parameterize $\sigma_0(x) = [1 + \exp(-\sum_i w_i \rho^i(x))]^{-1}$, then the eligibility of this actor parameter takes the form:

$$\frac{\partial}{\partial w_i} \ln \mu_{\theta, w}(x, u) = ((u - m_\theta(x))^2 - \sigma_0^2(x)) (1 - \sigma_0(x)) \rho^i(x) = \kappa_{\theta, w}^i(x, u) \rho^i(x).$$

It can be shown that this set of vectors forms a linearly independent basis set, which is also independent of the bases Ψ .

4 Candidate Features

In this section we investigate two different approaches for choosing linearly independent actor features, $\rho^j(\mathbf{x})$.

4.1 Barycentric Interpolation

Barycentric interpolants described in (Munos and Moore, 1998, 2002) are defined as an arbitrary set of (non-overlapping) mesh points ξ_i distributed across the state space. We denote the vector-valued output of the function approximator at each mesh point as $\mathbf{m}(\xi_i)$. For an arbitrary \mathbf{x} , if we define a simplex $S(\mathbf{x}) \in \{\xi_1, \dots, \xi_N\}$ such that \mathbf{x} is in the interior of the simplex, then the output at \mathbf{x} is given by interpolating the mesh points $\xi \in S(\mathbf{x})$:

$$\mathbf{m}_\theta(\mathbf{x}) = \sum_{\xi_i \in S(\mathbf{x})} \mathbf{m}(\xi_i) \lambda_{\xi_i}(\mathbf{x}).$$

Note that the interpolation is called *barycentric* if the positive coefficients $\lambda_{\xi_i}(\mathbf{x})$ sum to one, and if $\mathbf{x} = \sum_i \xi_i \lambda_{\xi_i}(\mathbf{x})$ (Munos and Moore, 1998). In addition, Munos and Moore (1998) denote the *piecewise linear* barycentric interpolation functions as functions for which the interpolation uses exactly $\dim(\mathbf{x}) + 1$ mesh points such that the simplex for state \mathbf{x} is the simplex which forms a triangulation of the state space and does not contain any interior mesh points. Barycentric interpolators are a popular representation for value functions, because they provide a natural mechanism for variable resolution discretization of the value function, and the barycentric co-ordinates allow the interpolators to be used directly by value iteration algorithms.

These interpolators also represent a linear function approximation architecture; we confirm here that the feature vectors are linearly independent. Let us consider the output at the mesh points as the parameters, $\theta_i = \mathbf{m}(\xi_i)$, and the interpolation function as the features $\rho_i(\mathbf{x}) = \lambda_{\xi_i}(\mathbf{x})$.

Proposition 3: The features $\rho_i(\mathbf{x})$ formed by the piecewise linear barycentric interpolation of a non-overlapping mesh ($\xi_i \neq \xi_j, \forall i \neq j$) form a linearly independent basis set.

Proof (sketch): For a non-overlapping mesh, consider the solution for the barycentric weights of a piecewise linear barycentric interpolation evaluated at $\mathbf{x} = \xi_i$. There are multiple simplices $S(\mathbf{x})$ that contain \mathbf{x} , but for each such simplex, \mathbf{x} is a vertex of that simplex. By definition of a simplex, \mathbf{x} is linearly independent of all other vertices of each simplex. As a result, the unique solution for the barycentric weights is $\rho_i(\mathbf{x}) = 1, \rho_j(\mathbf{x}) = 0, \forall j \neq i$. Since for each feature we can find an \mathbf{x} which is non-zero for only that feature, the basis set must be linearly independent.

Note that the traditional barycentric interpolators are *not* constrained to be linearly independent from the function 1.

4.2 Graph Laplacian

Proto-Value Functions (PVFs) (Mahadevan and Maggioni, 2006) have recently shown some success in automatic learning of representations in the context of function approximation in MDPs. In this approach, the agents learn global task-independent basis functions that reflect the large-scale geometry of the state-action space that all task-specific value functions must adhere to. Such basis functions are learned based on the topological structure of graphs representing the state (or state-action) space manifold. PVFs are essentially a subset of eigenfunctions of the graph *Laplacian* computed from a random walk graph generated by the agent. We show here that if the proto-value functions are used instead to represent features of the actor, instead of the critic, then this representation satisfies our Proposition 2.

Proposition 4: If the functions $\{\rho^j\}_{j=1}^n$ are the proto-value functions computed from the graph generated by a random walk in state space, then the set of critic features and the function 1 will form a linearly independent basis set, and will satisfy the weak form of the non-zero projection property presented in Equation 13.

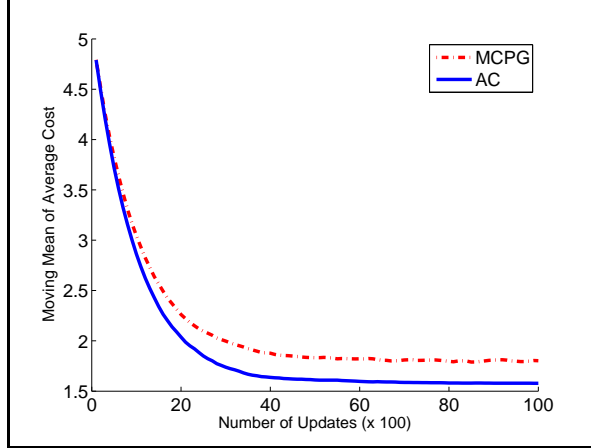


Figure 1: Performance comparison between MCPG and AC algorithm. These results are averaged over 5 trials. Note that this is an infinite-horizon average cost setting; the entire x-axis represents 1000 seconds of simulated data.

Proof (sketch): Since the functions $\{\rho^j\}_{j=1}^n$ are the eigen-functions of the graph Laplacian computed from the graph generated by a random walk in state space, they are linearly independent. Note that the function $\underline{1}$ is always the eigenfunction of any graph Laplacian associated with the eigenvalue 0. That implies the functions $\{\rho^j\}_{j=1}^n$ are also linearly independent of the function $\underline{1}$. Based on the results of the Proposition 2, the critic feature functions ϕ^{ij} are also linearly independent and also linearly independent of the function $\underline{1}$, and the features will satisfy the weak non-zero projection property.

5 Experiments

We demonstrate the effectiveness of our feature selections by learning a control policy for the swing-up task on a torque- limited inverted pendulum, governed by $\ddot{q} = \frac{1}{ml^2} [\tau - b\dot{q} - mgl \cos q]$, with $m = 1, l = 1, b = 1, g = 9.8, |\tau| < 3$, and initial conditions $q = -\frac{\pi}{2}, \dot{q} = 0$. We use an infinite-horizon, average reward formulation (no resetting) with the instantaneous cost function $g(q, \dot{q}, \tau) = \frac{1}{2}(q - \frac{\pi}{2})^2 + \frac{1}{20}\dot{q}^2 + \frac{1}{10}\tau^2$. The policy is evaluated every $dt = 0.1$ seconds; τ is held constant (zero-order hold) between evaluations.

Samples for our graph Laplacian are generated using rapidly-exploring randomized trees (LaValle and Kuffner, 2000) for coverage. Note that this is in place of the traditional “behavioral policy” used to identify the proto-value functions; it provides a fast and efficient coverage of our continuous state space.

For comparison, we also implemented the Markov Chain Policy Gradient algorithm (MCPG) (Algorithm 1 in (Baxter and Bartlett, 2001)). For both methods the policy is parametrized as in Equation 6 using PVFs as the basis set in the actor (i.e., the functions ρ^j). Figure 1 shows the moving mean of the average cost of the AC and MCPG algorithms, each averaged over five trials. Each trial starts the pendulum from the initial condition, with the parameters of the actor and critic initialized to small random values. We use a continuous setting in the AC experiment that consists of 10,000 steps (updates). In the MCPG experiment, each trial consists of 100 episodes of length 100 steps. At the beginning of each episode the pendulum is reset to the initial condition and accumulates the updates until the end of the episode where the policy parameters are adjusted using the updates collected throughout an episode. The key observation in this figure is that the AC method, using the policy structure that we described in Equation 6 together with PVFs, converges smoothly to a local minimum. The comparison with MCPG supports the promise of AC algorithms to outperform pure policy gradient methods. We also conducted an experiment where we used a polynomial basis set (instead of PVFs) which quickly led to a degenerate case.

6 Discussion

Our goal in this work has been to provide a mechanism for designing a minimal set of critic features which satisfy the conditions required for the convergence proof provided by Konda and Tsitsiklis (2003). It should be noted that meeting these conditions in a strict sense may not be necessary for obtaining stable and efficient actor-critic algorithms—these conditions were used to facilitate the convergence proofs. Independent of the convergence proof, a linearly independent basis set should generally permit faster learning than a linearly dependent basis set. Additionally, linear independence from the function $\underline{1}$ takes advantage of the well-known property that value functions can be offset by a constant value and retain the same gradient information for the policy; there is no policy gradient information in the value function along the direction of $\underline{1}$. Therefore, the qualities we have investigated appear to have general merit for representations used in value learning.

One could imagine other metrics which describe good critic features for actor-critic algorithms. For example, the graph Laplacians attempt to capture some of the geometry of the problem in a sparse representation. One could also investigate the qualities that are most desirable for the actor representation (besides permitting a good critic feature set). Although omitting $\underline{1}$ from the graph Laplacian bases is attractive because the remaining bases satisfy the non-zero projection property, when we choose to do this, we are certainly restricting the policy class. The inability to express a constant bias in the policy may be an undesirable penalty for some problems. It is also worth noting that there are other forms of the actor-critic algorithm which do not depend on a well-formed gradient spanning the the full value function to guarantee convergence (e.g., Kimura and Kobayashi (1998)).

7 Conclusions

In this paper, we provide some insights for designing features for actor-critic algorithms with function approximation. For a limited policy class, we demonstrate that a linearly independent feature set in the actor permits a linearly independent feature set in the critic. This condition is satisfied by the piecewise linear barycentric interpolators, and by the features based on a graph Laplacian. When combined with an additional linear independence with the function $\underline{1}$, the critic features for any particular θ are uniformly bounded away from zero. This condition is satisfied by the graph Laplacian features. Finally, our experimental results demonstrate that our proposed representation smoothly and efficiently converges to a local minimum for a simulated inverted pendulum control task.

8 Acknowledgements

The authors would like to thank John Tsitsiklis for a very helpful discussion. This work was supported by the DARPA Learning Locomotion program (AFRL contact # FA8650-05-C-7262).

References

- Barto, A., Sutton, R., and Anderson, C. (1983). Neuron-like elements that can solve difficult learning control problems. In *IEEE Trans. on Systems, Man and Cybernetics*, volume 13, pages 835–846.
- Baxter, J. and Bartlett, P. (2001). Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15:319–350.
- Kimura, H. and Kobayashi, S. (1998). An analysis of actor/critic algorithms using eligibility traces: Reinforcement learning with imperfect value functions. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 278–286.
- Konda, V. and Tsitsiklis, J. (2000). Actor-critic algorithms. In *Advances in Neural Information Processing Systems 12*.

- Konda, V. and Tsitsiklis, J. (2003). Actor-critic algorithms. *SIAM Journal on Control and Optimization*, 42(4):1143–1166.
- LaValle, S. and Kuffner, J. (2000). Rapidly-exploring random trees: Progress and prospects. In *Proceedings of the Workshop on the Algorithmic Foundations of Robotics*.
- Mahadevan, S. and Maggioni, M. (2006). Proto-value functions: A laplacian framework for learning representation and control in markov decision processes. Technical Report TR-2006-35, University of Massachusetts, Department of Computer Science.
- Marbach, P. and Tsitsiklis, J. (1998). Simulation-based optimization of markov reward processes.
- Munos, R. and Moore, A. (1998). Barycentric interpolators for continuous space and time reinforcement learning. In Kearns, M. S., Solla, S. A., and Cohn, D. A., editors, *Advances in Neural Information Processing Systems*, volume 11, pages 1024–1030. NIPS, MIT Press.
- Munos, R. and Moore, A. (2002). Variable resolution discretization in optimal control. *Machine Learning*, 49(2/3):291–323.
- Sutton, R., McAllester, D., Singh, S., and Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, pages 1057–1063.

